

Non-Asymptotic Delay Bounds for (k, l) Fork-Join Systems and Multi-Stage Fork-Join Networks

Markus Fidler
Institute of Communications Technology
Leibniz Universität Hannover

Yuming Jiang
Department of Telematics
NTNU Trondheim

Abstract—Parallel systems have received increasing attention with numerous recent applications such as fork-join systems, load-balancing, and l -out-of- k redundancy. Common to these systems is a join or resequencing stage, where tasks that have finished service may have to wait for the completion of other tasks so that they leave the system in a predefined order. These synchronization constraints make the analysis of parallel systems challenging and few explicit results are known. In this work, we model parallel systems using a max-plus approach that enables us to derive statistical bounds of waiting and sojourn times. Taking advantage of max-plus system theory, we also show end-to-end delay bounds for multi-stage fork-join networks. We contribute solutions for basic G|G|1 fork-join systems, parallel systems with load-balancing, as well as general (k, l) fork-join systems with redundancy. Our results provide insights into the respective advantages of l -out-of- k redundancy vs. load-balancing.

I. INTRODUCTION

Fork-join systems are an essential model of parallel data processing, such as Hadoop MapReduce [1], where jobs are divided into k tasks (fork) that are processed in parallel by k servers. Once all tasks of a job are completed, the results are combined (join) and the job leaves the system. Multi-stage fork-join networks comprise several fork-join systems in tandem, where all tasks of a job have to be completed at the current stage before the job is handed over to the next stage. The difficulty in analyzing such systems is due to a) the statistical dependence of the workload of the parallel servers that is due to the common arrival process [2], [3], and b) the synchronization that is enforced by the join operation [2].

Significant research has been performed to analyze the performance of fork-join systems. Exact results are, however, known for few specific systems only, such as for two parallel M|M|1 queues [4], [5]. For more complex systems, approximation techniques, e.g., [3], [5]–[11], and bounds, using stochastic orderings [2], martingales [12], or stochastic burstiness constraints [13], have been explored. Given the difficulties posed by single-stage fork-join systems, few works consider multi-stage networks. A notable exception is [10] where an approximation for closed fork-join networks is developed.

Related synchronization problems occur also in case of load balancing using parallel servers and in case of multi-path routing of packet data streams [14] using multi-path protocols [12]. The tail behavior of delays in multi-path routing is investigated

in [14] as well as in [15], [16] where large deviation results of resequencing delays for parallel M|M|1 servers are derived.

A further synchronization constraint applies in split-merge systems, that are a variant of fork-join systems, where all tasks of a job have to start execution at the same time. In contrast, in a fork-join system, the start time of tasks is not synchronized. Split-merge systems are solvable to some extent as they can be expressed as a single server queue, where the service process is governed by the maximal service time [6], [12], [17], [18].

A generalization of fork-join systems are (k, l) fork-join systems, where a job is finished once any l out of its k tasks are completed. The model enables analyzing systems with redundant servers and suitable coding. A first queueing model of maximum distance separable (MDS) codes is presented in [19] and the download of coded content from distributed storage systems is modelled as a (k, l) fork-join system in [18]. Generally, (k, l) fork-join systems are governed by the l th order statistic. The authors of [18] compute bounds of the mean response time using as approximation a (k, l) split-merge model with Poisson arrivals, that is solved by the Pollaczek-Khinchin formula.

Most closely related to this work are the two recent papers [12], [13] that employ similar methods. The work [13] considers single-stage fork-join systems with load balancing, general arrivals of the type defined in [20], and deterministic service. A service curve characterization of fork-join systems is provided and first statistical delay bounds are presented. The paper [12] contributes delay bounds for single-stage fork-join systems with renewal as well as Markov modulated inter-arrival times and independent and identically distributed (iid) service times. The authors prove that delays for fork-join systems grow as $\mathcal{O}(\ln k)$ for k parallel servers as also found in [2]. Split-merge systems exhibit an inferior performance, where the stability region of k parallel M|M|1 servers is shown to decrease with $\ln k$. The work also includes a first application to multi-path routing, assuming a simple window-based protocol that operates on batches of packets. Response times are computed for entire batches and the authors conclude that multi-path routing is beneficial in case of only two parallel paths and moderate to high utilization. Otherwise, resequencing delays are found to dominate.

While [12] focuses on split-merge vs. fork-join systems with iid service times, we consider also the case of non-iid service, where we are able to generalize important results, such as

This work was supported in part by the European Research Council (ERC) under Starting Grant UniQue (StG 306644).

the growth of delays in $\mathcal{O}(\ln k)$ for fork-join systems with k parallel servers. Further, we investigate advanced fork-join systems beyond [12], [13], including (k, l) -fork join systems and multi-stage fork-join networks, where we present a scaling of end-to-end delay bounds in $\mathcal{O}(h \ln(hk))$ for h stages. Considering heterogeneous servers and deterministic as well as random thinning of arrival processes, we give insights into load-balancing. We draw essential conclusions on the advantages of load-balancing and l -out-of- k redundancy.

The remainder of this paper is structured as follows. In Sec. II, we phrase basic models of $G|G|1$ as well as $GI|GI|1$ fork-join systems in max-plus system theory and show an application to load balancing. We extend the model to parallel servers with thinning and resequencing, as it applies, e.g., in multi-path routing, in Sec. III. In Sec. IV, we consider advanced (k, l) fork-join systems and multi-stage fork-join networks. Sec. V presents brief conclusions.

II. BASIC FORK-JOIN SYSTEMS

In this section, we derive a set of results for basic fork-join systems. Compared to [12], we define a general queueing model in max-plus system theory [21]–[25] that is a branch of the deterministic [22], [26], [27], respectively, stochastic network calculus [22], [28]–[34]. The model enables us to express also more advanced fork-join systems and networks in the following sections. Further, we generalize central results from [12] to the case of $G|G|1$ queues. For the special case of $GI|GI|1$ queues we recover the bounds from [12].

A. Fork-Join System Model

We label jobs by n where $n \geq 1$ and denote $A(n)$ the time of arrival of job n . For notational convenience, we define $A(0) = 0$. Further, we let $A(m, n) = A(n) - A(m)$ be the time between the arrival of job m and job n for $n \geq m \geq 1$. It follows that $A(\nu, \nu + 1)$ denotes the inter-arrival time between job ν and job $\nu + 1$ for $\nu \geq 1$, so that we can also express

$$A(m, n) = \sum_{\nu=m}^{n-1} A(\nu, \nu + 1) \quad (1)$$

as a sum of inter-arrival times. By definition, $A(n, n) = 0$ for $n \geq 1$. Similarly, $D(n)$ defines the departure time of job n .

We denote $S(n)$ the service time of job n for $n \geq 1$. The cumulative service time for jobs m up to and including n follows for $n \geq m \geq 1$ as

$$S(m, n) = \sum_{\nu=m}^n S(\nu). \quad (2)$$

For a lossless, work-conserving, first-in first-out (fifo) system, the departures for all $n \geq 1$ are [22]

$$D(n) = \max_{\nu \in [1, n]} \{A(\nu) + S(\nu, n)\}. \quad (3)$$

To see this, consider all jobs $[\nu, n]$ that belong to the same busy period as job n , i.e., job ν arrives at time $A(\nu)$ at an empty system that is continuously busy afterwards. It follows that $D(n) = A(\nu) + S(\nu, n)$ since starting at $A(\nu)$ all jobs $[\nu, n]$

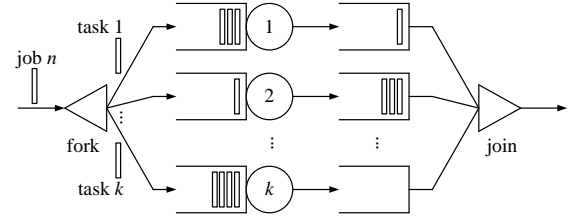


Fig. 1. Fork-join system. Each job is composed of k tasks with individual service requirements, that are mapped to k fifo servers (fork). Once all tasks of a job are completed, the job leaves the system (join), i.e., the tasks of a job wait at the join stage until all tasks of the job are completed.

have to be served until job n departs, requiring a cumulative service time of $S(\nu, n)$. In general ν is a priori unknown, however, it is known that there exists $\nu \geq 1$ such that $D(n) = A(\nu) + S(\nu, n)$ for $n \geq 1$ so that it can be concluded that

$$D(n) \leq \max_{\nu \in [1, n]} \{A(\nu) + S(\nu, n)\}. \quad (4)$$

Further, $D(n) \geq D(\nu) + S(\nu + 1, n)$ for all $n > \nu \geq 1$. Since $D(\nu) \geq A(\nu) + S(\nu)$ it follows that $D(n) \geq A(\nu) + S(\nu, n)$ for all $n \geq \nu \geq 1$. Consequently, $D(n) \geq \max_{\nu \in [1, n]} \{A(\nu) + S(\nu, n)\}$ so that combined with (4) we obtain (3). For many applications the upper bound (4) is sufficient as it will provide upper bounds of waiting and sojourn times.

For the sojourn time of job n defined as $T(n) = D(n) - A(n)$ it follows by insertion of (3) for $n \geq 1$ that

$$T(n) = \max_{\nu \in [1, n]} \{S(\nu, n) - A(\nu, n)\}. \quad (5)$$

For the waiting time, defined as $W(n) = [D(n-1) - A(n)]^+$ for $n \geq 1$, where $[X]^+ = \max\{X, 0\}$ is the non-negative part and $D(0) = 0$ by definition, it holds for $n \geq 1$ that

$$W(n) = \left[\sup_{\nu \in [1, n-1]} \{S(\nu, n-1) - A(\nu, n)\} \right]^+. \quad (6)$$

Here, we use sup since for $n = 1$ (6) evaluates an empty set. For non-negative real numbers the sup of an empty set is zero.

In a fork-join system, see Fig. 1, each job $n \geq 1$ is composed of k tasks with service times $S_i(n)$ for $i \in [1, k]$, i.e., the service requirements of the tasks may differ from each other. The tasks are mapped to k parallel servers (fork) and once all tasks are served, the job leaves the system (join). The parallel servers are not synchronized, i.e., server $i \in [1, k]$ starts serving task i of job $n+1$ if any, once it finished serving task i of job n that departs from server i at $D_i(n)$. A job n is finished once all of its tasks $i \in [1, k]$ are finished, i.e.,

$$D(n) = \max_{i \in [1, k]} \{D_i(n)\}.$$

Hence, the sojourn time is $T(n) = \max_{i \in [1, k]} \{D_i(n) - A(n)\}$ for $n \geq 1$ and by insertion of (3) for each server $i \in [1, k]$

$$T(n) = \max_{i \in [1, k]} \left\{ \max_{\nu \in [1, n]} \{S_i(\nu, n) - A(\nu, n)\} \right\}. \quad (7)$$

Similarly, a waiting time that considers the task that starts service last can be defined as a $\max_{i \in [1, k]}$ of (6).

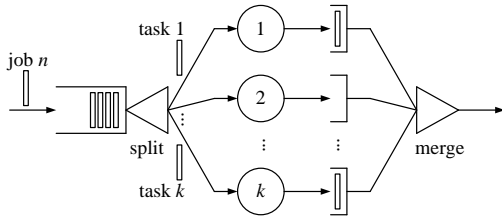


Fig. 2. Split-merge system. Compared to the fork-join system, tasks have an additional synchronization constraint, i.e., the execution of the tasks of a job has to start at the same time.

In a split-merge system on the other hand, see Fig. 2, the execution of all tasks of a job starts simultaneously, i.e., once server i finishes task i of job n it idles until all tasks $j \in [1, k]$ of that job are finished, before execution of the tasks of job $n + 1$ if any starts. The combined service and idle times at each of the servers follows as $\max_{i \in [1, k]} \{S_i(n)\}$ such that the split-merge system as a whole can be expressed as a single server system that is governed by (3) with service process

$$S(m, n) = \sum_{\nu=m}^n \max_{i \in [1, k]} \{S_i(\nu)\}. \quad (8)$$

The general problem of split-merge systems is that $\max_{i \in [1, k]} \{S_i(\nu)\}$ increases with k , with few exceptions such as in case of identical task service times. The increase implies longer idle times that result in a reduced stability region.

B. Performance Bounds for GI|GI|1 and G|G|1 servers

Next, we derive statistical performance bounds for the fork-join systems defined above. We consider the general model of G|G|1 servers. The results enable us to generalize recent findings obtained for iid service times, i.e., for a GI service model, in [12]. Like [12], we assume that the inter-arrival times of jobs are independent of their service times. We do, however, not require that the service times of the tasks of a job are independent. We will show how to benefit from statistically independent tasks in Sec. IV-A.

We consider arrival and service processes that belong to the broad class of (σ, ρ) -constrained processes [22], that are characterized by affine bounding functions of the moment generating function (MGF). The MGF of a random variable X is defined as $M_X(\theta) = \mathbb{E}[e^{\theta X}]$ where θ is a free parameter. The following definition adapts [22] to max-plus systems.

Definition 1. An arrival process is (σ_A, ρ_A) -lower constrained if for all $n \geq m \geq 1$ and $\theta > 0$ it holds that

$$\mathbb{E}[e^{-\theta A(m, n)}] \leq e^{-\theta(\rho_A(-\theta)(n-m) - \sigma_A(-\theta))}.$$

Similarly, a service process is (σ_S, ρ_S) -upper constrained if for all $n \geq m \geq 1$ and $\theta > 0$ it holds that

$$\mathbb{E}[e^{\theta S(m, n)}] \leq e^{\theta(\rho_S(\theta)(n-m+1) + \sigma_S(\theta))}.$$

For the special case of iid inter-arrival and service times, i.e., in case of GI|GI|1 servers, we have $\mathbb{E}[e^{-\theta A(n, n+1)}] = \mathbb{E}[e^{-\theta A(1, 2)}]$ and $\mathbb{E}[e^{\theta S(n)}] = \mathbb{E}[e^{\theta S(1)}]$ for all $n \geq 1$. Further,

since the MGF of a sum of independent random variables is the product of their MGFs, i.e., $M_{X+Y}(\theta) = M_X(\theta)M_Y(\theta)$, minimal traffic and service parameters can be derived from (1) and (2) as $\sigma_A(-\theta) = 0$, $\sigma_S(\theta) = 0$,

$$\rho_A(-\theta) = -\frac{1}{\theta} \ln \mathbb{E}[e^{-\theta A(1, 2)}], \quad (9)$$

and

$$\rho_S(\theta) = \frac{1}{\theta} \ln \mathbb{E}[e^{\theta S(1)}]. \quad (10)$$

Parameter $\rho_A(-\theta)$ decreases with $\theta > 0$ from the mean to the minimum inter-arrival time and $\rho_S(\theta)$ increases from the mean to the maximum service time.

Theorem 1. Consider a fork-join system with $i \in [1, k]$ parallel servers, each with arrivals and service as specified by Def. 1. It holds for the waiting time for all $n \geq 1$ that

$$\mathbb{P}[W(n) > \tau] \leq \sum_{i=1}^k \alpha_i e^{-\theta_i \tau},$$

and for the sojourn time that

$$\mathbb{P}[T(n) > \tau] \leq \sum_{i=1}^k \alpha_i e^{\theta_i \rho_{S_i}(\theta_i)} e^{-\theta_i \tau}.$$

In the general case of G|G|1 servers, the free parameters $\theta_i > 0$ have to satisfy $\rho_{S_i}(\theta_i) < \rho_A(-\theta_i)$ for $i \in [1, k]$ and

$$\alpha_i = \frac{e^{\theta_i(\sigma_A(-\theta_i) + \sigma_{S_i}(\theta_i))}}{1 - e^{-\theta_i(\rho_A(-\theta_i) - \rho_{S_i}(\theta_i))}}.$$

In the special case of GI|GI|1 servers, the $\theta_i > 0$ have to satisfy $\rho_{S_i}(\theta_i) \leq \rho_A(-\theta_i)$ and $\alpha_i = 1$ for $i \in [1, k]$.

For the special case of GI|GI|1 servers, Th. 1 recovers the bounds from [12]. Further, for $k = 1$ the classical bound for the waiting time of a single GI|GI|1 server [35] is obtained in the max-plus system theory. Like [35], the proof uses Doob's martingale inequality [36]. The proof for the G|G|1 server brings the approach from [22], [30] forward to max-plus fork-join systems. The important property of the G|G|1 result is that it differs only by constants α_i from the GI|GI|1 result and otherwise recovers the characteristic exponential tail decay $e^{-\theta_i \tau}$ with decay rate θ_i . We note that Th. 1 does not make an assumption of independence regarding the parallel servers. Indeed, independence cannot be assumed as the waiting and sojourn times of the individual servers depend on the common arrival process [2], [3]. The proof is provided in the appendix.

To obtain a solution for split-merge systems, recall that the system can be expressed as a single server with service process (8). A corresponding service constraint as in Def. 1 can be directly inserted into Th. 1 where $k = 1$. For iid service times it can be seen from (8) and (10) that

$$\rho_S(\theta) = \frac{1}{\theta} \ln \mathbb{E}[e^{\theta \max_{i \in [1, k]} \{S_i(1)\}}] \leq \frac{1}{\theta} \ln \left(\sum_{i=1}^k \mathbb{E}[e^{\theta S_i(1)}] \right).$$

has at most a logarithmic growth with the number of parallel servers, resulting in a corresponding reduction of the stability

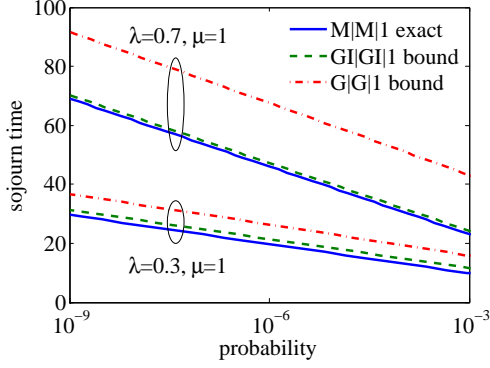


Fig. 3. M|M|1 server. The bounds show the correct exponential tail decay.

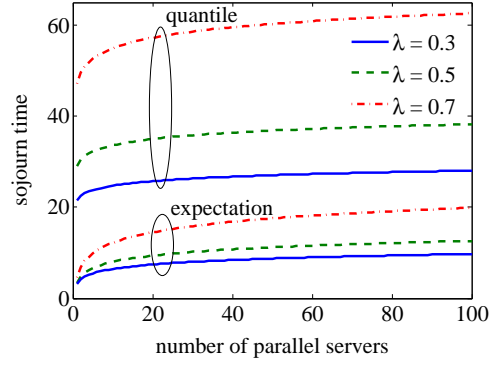


Fig. 4. Fork-join system. Sojourn time bounds grow with $\ln k$ for k servers.

region. A decrease of the stability region with $\ln k$ is also shown in [12], where the authors advise against split-merge implementations based on an in-depth comparison with fork-join systems. In the sequel, we will focus on fork-join systems.

To investigate the scaling of fork-join systems, we consider the homogeneous case where the service times $S_i(\nu)$ for $i \in [1, k]$ are identically distributed with parameters $\rho_{S_i} = \rho_S$, but not necessarily independent. From Th. 1 it follows for the sojourn time that

$$\mathbb{P}[T(n) > \tau] \leq k\alpha e^{\theta \rho_S(\theta)} e^{-\theta \tau}. \quad (11)$$

After some reordering, the bound of the sojourn time

$$\mathbb{P}\left[T(n) > \tau + \frac{\ln k}{\theta}\right] \leq \alpha e^{\theta \rho_S(\theta)} e^{-\theta \tau}$$

shows a growth with $\ln k$. By integration of the tail of (11), where we use that $\mathbb{P}[T(n) > \tau] \leq 1$,

$$\begin{aligned} \mathbb{E}[T(n)] &= \int_0^\infty \mathbb{P}[T(n) > \tau] d\tau \\ &\leq \int_0^{\tau^*} d\tau + k\alpha e^{\theta \rho_S(\theta)} \int_{\tau^*}^\infty e^{-\theta \tau} d\tau, \end{aligned}$$

and $\tau^* = \ln(k\alpha e^{\theta \rho_S(\theta)})/\theta$, the expected sojourn time

$$\mathbb{E}[T(n)] \leq \rho_S(\theta) + \frac{\ln(k\alpha) + 1}{\theta} \quad (12)$$

is also limited by $\ln k$. The result applies for general arrival and service processes and generalizes the finding of $\ln k$ that is obtained in [2], [12] for iid service times. The growth is larger for smaller θ corresponding to a higher utilization.

M|M|I Servers: For illustration, we consider iid exponential inter-arrival and service times with parameters λ and μ , respectively. The normalized log-MGFs (9) and (10) are

$$\rho_A(-\theta) = -\frac{1}{\theta} \ln\left(\frac{\lambda}{\lambda + \theta}\right) \text{ and } \rho_S(\theta) = \frac{1}{\theta} \ln\left(\frac{\mu}{\mu - \theta}\right), \quad (13)$$

where $\theta \in (0, \mu)$. From the condition $\rho_S(\theta) \leq \rho_A(-\theta)$ it follows that $\theta \leq \mu - \lambda$ under the stability condition $\mu > \lambda$. By choice of the maximal $\theta = \mu - \lambda$ we have from (11) that

$$\mathbb{P}[T(n) > \tau] \leq k \frac{\mu}{\lambda} e^{-(\mu - \lambda)\tau}. \quad (14)$$

First, we consider the accuracy of the bounds from Th. 1 for $k = 1$ that is the case of a single M|M|1 server. Compared to the exact distribution of the sojourn time of the M|M|1 server, that is $\mathbb{P}[T(n) > \tau] \leq e^{-(\mu - \lambda)\tau}$ see, e.g., [37], (14) has the same tail decay, but shows a different pre-factor, i.e., μ/λ . Obviously, the bound becomes better if the utilization is high so that μ/λ approaches one. In Fig. 3, we illustrate the bounds from Th. 1 for a single server compared to the exact M|M|1 result. Clearly, the curves show the same tail decay, where the GI|GI|1 bound provides better numerical accuracy compared to the G|G|1 bound that does not use independence of the increment processes and hence has parameter $\alpha > 1$.

In Fig. 4, we consider $k \geq 1$ parallel servers with $\mu_i = 1$ and show bounds of the expected sojourn time and sojourn time quantiles τ , where $\mathbb{P}[T(n) \geq \tau] \leq \varepsilon$ and $\varepsilon = 10^{-6}$. The curves show the characteristic logarithmic growth with k .

For reasons of space, we will frequently employ the M|M|1 model (13) for numerical evaluations. We note that results of the same type can be derived from Th. 1 for G|G|1 servers with parameters specified by Def. 1.

C. Load Balancing

Next, we consider heterogeneous GI|GI|1 servers, i.e., the tasks of a job can have service requirements $S_i(n)$ for $i \in [1, k]$ with different parameters. Also, the parallel servers can have different capacities c_i . The resulting service times become $S_i(n)/c_i$ and the MGF follows for $n \geq 1$ as

$$\mathbb{M}_{S_i(n)/c_i}(\theta_i) = \mathbb{E}[e^{\theta_i S_i(1)/c_i}] = \mathbb{M}_{S_i(1)}(\theta_i/c_i). \quad (15)$$

By application of Th. 1 it holds for all $n \geq 1$ that

$$\mathbb{P}[T(n) > \tau] \leq \sum_{i=1}^k e^{\theta_i \rho_{S_i/c_i}(\theta_i)} e^{-\theta_i \tau}, \quad (16)$$

for all $\theta_i > 0$ that satisfy $\rho_{S_i/c_i}(\theta_i) \leq \rho_A(-\theta_i)$.

In the following we consider two different strategies to assign the capacities c_i for load balancing. For ease of exposition, we assume that the c_i are continuous.

- 1) Assign capacity proportional to the mean service time requirements of the tasks. As a consequence, the mean service times and the utilization of each of the parallel

servers becomes identical. The allocation is independent of the inter-arrival times. It is stable as long as the mean inter-arrival time is larger than the mean service time.

- 2) Assign capacity such that the same statistical bound of the sojourn time is achieved for all servers. This allocation may differ from the first one, so that the parallel servers may not have identical utilization.

M|M|I Servers: We start with the basic case of exponential inter-arrival and service times with parameters λ and μ_i , respectively. In case of strategy 1, capacity is allocated proportionally to the mean service time requirements, i.e., $c_i = v/\mu_i$ where v is a positive constant. Given the sum of the capacities is limited as $c = \sum_{i=1}^k c_i$, we have $v = c/\sum_{i=1}^k 1/\mu_i$. By insertion of (15) into (10) it follows with (13) that

$$\rho_{S_i/c_i}(\theta_i) = \frac{1}{\theta_i} \ln \left(\frac{\mu_i}{\mu_i - \theta_i/c_i} \right) = \frac{1}{\theta_i} \ln \left(\frac{v}{v - \theta_i} \right).$$

To ensure $\rho_{S_i/c_i}(\theta_i) \leq \rho_A(-\theta_i)$, where $\rho_A(-\theta_i)$ is given by (13), it has to hold that $\theta_i \leq v - \lambda$ under the stability condition that $v > \lambda$. By choice of the maximal $\theta_i = v - \lambda$, (16) evaluates to

$$P[T(n) > \tau] \leq k \frac{v}{\lambda} e^{-(v-\lambda)\tau}.$$

The result shows that strategy 1 in case of M|M|I servers also achieves the target of strategy 2, i.e., it implements identical bounds of the sojourn time for all k servers.

Gaussian Servers: Secondly, we consider Gaussian inter-arrival and service times with mean and variance η_A, ς_A^2 and η_S, ς_S^2 , respectively. The normalized log-MGFs (9) and (10) are

$$\rho_A(-\theta) = \eta_A - \frac{\theta}{2} \varsigma_A^2 \quad \text{and} \quad \rho_S(\theta) = \eta_S + \frac{\theta}{2} \varsigma_S^2.$$

Considering heterogeneous service parameters $\eta_{S_i}, \varsigma_{S_i}^2$, strategy 1 implements a capacity allocation $c_i = v\eta_{S_i}$, where v is a positive constant as above. With (15) we obtain

$$\rho_{S_i/c_i}(\theta_i) = \frac{\eta_{S_i}}{c_i} + \frac{\theta_i}{2} \left(\frac{\varsigma_{S_i}}{c_i} \right)^2 = \frac{1}{v} + \frac{\theta_i}{2} \left(\frac{\varsigma_{S_i}}{v\eta_{S_i}} \right)^2.$$

A bound of the sojourn time follows from (16) where we choose the maximal θ_i that satisfy $\rho_{S_i/c_i}(\theta_i) \leq \rho_A(-\theta_i)$ for all $i \in [1, k]$. Clearly, as the parameters $\rho_{S_i/c_i}(\theta_i)$ are heterogeneous, the maximal tail decay parameters θ_i will in general not be identical. As a consequence, the capacity allocation does in this case not achieve the goal of strategy 2.

To implement strategy 2, (16) implies that the capacity is allocated in a way such that identical parameters $\rho_{S_i/c_i}(\theta_i)$ are achieved for the same target tail decay parameter $\theta_i = \theta$ for all servers $i \in [1, k]$. Considering the minimal capacity allocation that satisfies $\rho_{S_i/c_i}(\theta) \leq \rho_A(-\theta)$, the capacities c_i follow readily as the solution of the quadratic problem

$$\frac{\eta_{S_i}}{c_i} + \frac{\theta}{2} \left(\frac{\varsigma_{S_i}}{c_i} \right)^2 = \eta_A - \frac{\theta}{2} \varsigma_A^2.$$

The important observation is that the obvious strategy 1, that balances the average utilization, does in general not achieve

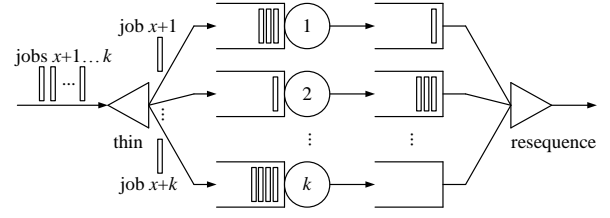


Fig. 5. System with thinning and resequencing. Compared to a fork-join system, jobs are not composed of tasks. Instead, entire tasks are mapped to the parallel servers, e.g., deterministically in round robin order.

the same bound of the sojourn time for all servers, i.e., under strategy 1 certain servers may be late frequently.

III. THINNING AND RESEQUENCING

We investigate systems of k parallel servers with thinning and resequencing, as it applies, e.g., in case of multi-path routing. The difference to fork-join systems is that jobs are not composed of tasks that are served in parallel but instead each job is mapped in its entirety to one of the parallel servers. As a consequence, the external arrival process $A(n)$ is divided into k thinned processes $A_i(m)$. The numbering of jobs is such that $A(n)$ is the arrival time of job n before thinning, whereas $A_i(m)$ denotes the arrival time of the m th job of the thinned process at system i with service time $S_i(m)$. The departures $D_i(m)$ are resequenced in the original order of $A(n)$ to form the departure process $D(n)$. We contribute solutions for random as well as deterministic thinning. An example of a system with deterministic thinning is depicted in Fig. 5

A. Random Thinning

In case of random thinning, each job is mapped to one of the k servers according to iid discrete (not necessarily uniform) random variables with support $[1, k]$. From the iid property it follows that the mapping of each job to a certain server $i \in [1, k]$ is an independent Bernoulli trial with parameter p_i . Denote $X_i(m)$ the number of the job that becomes the m th job that is mapped to server i . It follows that $X_i(m)$ is a sum of m iid geometric random variables with parameter p_i , i.e., $X_i(m)$ is negative binomial. The arrival process of server i is

$$A_i(m) = A(X_i(m)), \quad (17)$$

for $m \geq 1$. Conversely, given jobs $[1, n]$ of the external arrival process, denote $Y_i(n)$ the quantity of jobs that are mapped to server i . It follows that $Y_i(n)$ is binomial with parameter p_i . Considering fifo servers with arrival processes (17) and departure processes $D_i(m)$ for $i \in [1, k]$, the combined in-sequence departure process for $n \geq 0$ follows as

$$D(n) = \max_{i \in [1, k]} \{D_i(Y_i(n))\}, \quad (18)$$

where $D_i(0) = 0$ by convention. Note that (18) has to verify only the departure of job $Y_i(n)$ of each server $i \in [1, k]$ since the departure of job $Y_i(n)$ from server i implies the departure of all jobs $\nu \in [1, Y_i(n)]$ of the same server due to fifo order.

Next, we substitute (3) for $D_i(Y_i(n))$ in (18) and use (17) to derive the sojourn time $T(n) = D(n) - A(n)$ for $n \geq 1$ as

$$T(n) = \max_{i \in [1, k]} \left\{ \sup_{\nu \in [1, Y_i(n)]} \{S_i(\nu, Y_i(n)) - A(X_i(\nu), n)\} \right\},$$

where $\sup\{\emptyset\} = 0$. Further, since $A(n) \geq A_i(Y_i(n))$ for all $i \in [1, k]$ and $n \geq 1$ and with (17), we can estimate

$$T(n) \leq \max_{i \in [1, k]} \left\{ \sup_{\nu \in [1, Y_i(n)]} \{S_i(\nu, Y_i(n)) - A_i(\nu, Y_i(n))\} \right\}. \quad (19)$$

Due to the similarity of (19) with (7), the derivation of delay bounds closely follows the proof of Th. 1 and recovers the same result with one essential difference: Instead of the arrival process $A(n)$ the thinned processes $A_i(m)$ have to be considered for each of the servers $i \in [1, k]$, i.e., parameters (σ_A, ρ_A) in Th. 1 are replaced by the parameters $(\sigma_{A_i}, \rho_{A_i})$ of the thinned processes. The main effect of thinning is an increase of the parameters θ_i that result in a faster tail decay.

The thinned arrivals (17) are expressed as doubly random processes. Considering iid inter-arrival times the MGF of the thinned process is

$$M_{A_i(\nu, \nu+1)}(-\theta) = E \left[(M_{A(1,2)}(-\theta))^{X_i(1)} \right],$$

for $\nu \geq 1$. It follows after some reordering that

$$M_{A_i(\nu, \nu+1)}(-\theta) = M_{X_i(1)}(\ln M_{A(1,2)}(-\theta)). \quad (20)$$

Since $X_i(1)$ is a geometric random variable with MGF $M_{X_i(1)}(\theta) = p_i e^\theta / (1 - (1 - p_i)e^\theta)$ for $\theta < -\ln(1 - p_i)$, we obtain by insertion of (20) into (9) that

$$\rho_{A_i}(-\theta) = -\frac{1}{\theta} \ln \left(\frac{p_i M_{A(1,2)}(-\theta)}{1 - (1 - p_i) M_{A(1,2)}(-\theta)} \right), \quad (21)$$

where $\theta > 0$ so that $M_{A(1,2)}(-\theta) < 1/(1 - p_i)$.

B. Deterministic Thinning

In case of deterministic thinning, the random processes $X_i(m)$ and $Y_i(n)$ are replaced by deterministic functions. A round robin assignment of the jobs of an arrival process $A(n)$ to k servers results in the processes $A_i(m)$ given by (17) where

$$X_i(m) = k(m - 1) + i, \quad (22)$$

for $m \geq 1$ and $i \in [1, k]$. The combined in-sequence departure process $D(n)$ follows from (18) with

$$Y_i(n) = \left\lceil \frac{n - i + 1}{k} \right\rceil, \quad (23)$$

for $n \geq 1$ and $i \in [1, k]$. To see this, note that job n of the external arrival process becomes the $m = \lceil n/k \rceil$ th job of server $j = (n - 1) \bmod k + 1$. Hence, $Y_i(n) = m$ for $i \leq j$ and $Y_i(n) = m - 1$ for $i > j$ due to the round robin procedure, as can be verified for (23).

Statistical sojourn time bounds follow as in case of random thinning from (19) as a variant of Th. 1, where the parameters $(\sigma_{A_i}, \rho_{A_i})$ of the thinned arrival processes are

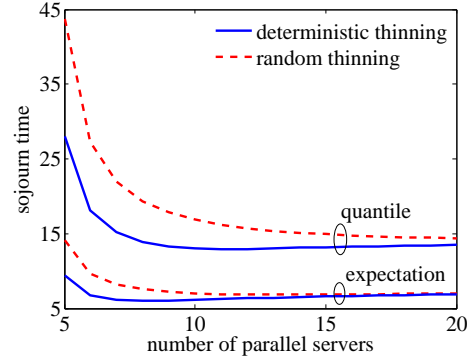


Fig. 6. Parallel servers achieve significant performance improvements if the utilization is high. Deterministic thinning outperforms random thinning.

used. Considering iid inter-arrival times, the MGFs of the deterministically thinned processes can be straightforwardly computed as $M_{A_i(\nu, \nu+1)}(-\theta) = (M_{A(1,2)}(-\theta))^k$ for $\nu \geq 1$. The parameters $\rho_{A_i}(-\theta)$ for $\theta > 0$ follow from (9) as

$$\rho_{A_i}(-\theta) = k \rho_A(-\theta). \quad (24)$$

For Markov arrival processes, an expression for deterministic thinning is given in [12].

$E_k[M/I]$ Servers: We consider arrivals $A(n)$ with iid exponential inter-arrival times with parameter λ . Deterministic thinning results in processes $A_i(n)$ where the inter-arrival times are a sum of k exponential random variables that is Erlang- k distributed. It follows from (24) for $\theta > 0$ that

$$\rho_{A_i}(-\theta) = -\frac{k}{\theta} \ln \left(\frac{\lambda}{\lambda + \theta} \right).$$

In case of random thinning the inter-arrival times of the thinned processes are exponentially distributed with parameter $p_i \lambda$. With $p_i = 1/k$, we have from (21) for $\theta > 0$ that

$$\rho_{A_i}(-\theta) = -\frac{1}{\theta} \ln \left(\frac{\lambda}{\lambda + k\theta} \right).$$

Given the service times at servers $i \in [1, k]$ are exponential with parameter $\mu_i = \mu$, we have $\rho_{S_i}(\theta) = \ln(\mu/(\mu - \theta))/\theta$ for $\theta \in (0, \mu)$. We choose the maximal parameter $\theta \in (0, \mu)$ such that $\rho_{S_i}(\theta) \leq \rho_{A_i}(-\theta)$ for all $i \in [1, k]$ and obtain bounds of the sojourn time from (11) and (12).

In Fig. 6, we compare deterministic and random thinning, where we show bounds of the expected sojourn time and of the sojourn time quantile with $\varepsilon = 10^{-3}$. The parameters are $\lambda = 4$ and $\mu = 1$, i.e., $k \geq 5$ parallel servers are required. Adding few more servers provides a significant advantage that is due to the decreasing utilization. The advantage diminishes, however, if k becomes large and eventually the delay bounds start to grow due to potential waiting in the resequencing stage. Note that Th. 1 does not assume independence of the parallel servers. Deterministic thinning generally performs better than random thinning, where the advantage is larger in case of smaller k .

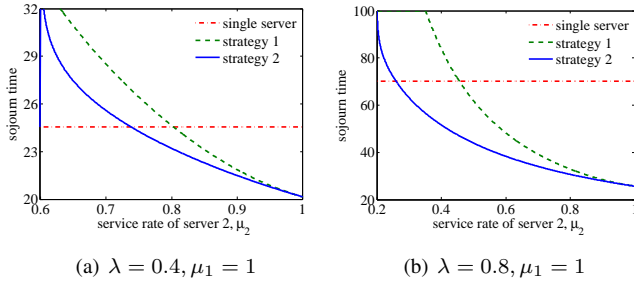


Fig. 7. The effectiveness of load balancing depends largely on the utilization.

C. Load Balancing

We consider load balancing for heterogeneous servers and investigate two basic strategies: Divide the arrivals so that

- 1) all servers have the same average utilization,
- 2) all servers have the same maximal tail decay.

The strategies are similar to those that we considered for fork-join systems in Sec. II-C. The difference is that we now thin the arrival process accordingly to achieve the goal.

M|M|I Servers: We consider the case of exponential inter-arrival times with parameter λ and exponential service times with parameters μ_i for servers $i \in [1, k]$. Random thinning is used to divide the arrivals into k sub-processes with parameters λ_i where $\sum_{i=1}^k \lambda_i = \lambda$. For strategy 1 the target is to achieve $\lambda_i/\mu_i = \lambda_j/\mu_j$ for all $i, j \in [1, k]$. It follows that $\lambda_i = \lambda\mu_i/(\sum_{j=1}^k \mu_j)$ for all $i \in [1, k]$. Strategy 2 seeks to equalize the maximal tail decay parameters $\theta_i = \mu_i - \lambda_i$, see (13), so that $\theta_i = \theta_j$ for all $i, j \in [1, k]$. A solution can be obtained iteratively as $\lambda_i = \mu_i - (\sum_{j=1}^k \mu_j - \lambda)/k$ for all $i \in [1, k]$. Servers that are assigned $\lambda_i < 0$ are excluded in the next iteration until all λ_i are non-negative.

Fig. 7 shows bounds of the sojourn time for $\varepsilon = 10^{-6}$. We consider $k = 2$ servers with parameters $\mu_1 = 1$, $\mu_2 \in [0.5, 1]$, and $\lambda = 0.4$ and 0.8 , respectively. A practical application of the model is, e.g., the multi-path transmission of a packet data stream over Wifi and cellular in parallel. For comparison the sojourn time for the case without load balancing, where only server one is used, is shown. Clearly, load balancing performs very well if the utilization is high. On the other hand, an advantage is achieved only if the capacity of server two is sufficiently large, i.e., in certain cases it is better to exclude slow servers. Strategy 2 performs better than strategy 1 as the allocation takes the tail decay of the sojourn time into account. If the servers are homogeneous, strategies 1 and 2 are identical.

IV. ADVANCED FORK-JOIN SYSTEMS AND NETWORKS

Next, we extend the methods from Sec. II to include (k, l) fork-join systems as well as multi-stage fork-join networks.

A. (k, l) Fork-Join Systems

An important generalization of fork-join systems are (k, l) fork-join systems, where a job is completed once l out of the k tasks of a job are finished. An application is the download of coded media from distributed storage systems [18]. Compared

to [18], we do not assume that the remaining $k - l$ unfinished tasks of a job are dropped once l tasks of the job are finished. This applies, e.g., in case of the transmission of encoded redundant data streams over multiple paths. Another practical example are systems with redundant jobs that correspond to $(k, 1)$ fork-join systems.

The advantage of (k, l) fork-join systems becomes evident if the service times at each of the k parallel servers are statistically independent. The result of Th. 1 does, however, not use an assumption of independence as the individual waiting and sojourn times of the parallel servers are stochastically dependent on the common arrival process [2], [3]. Hence, in order to take advantage of independent service times $S_i(n)$ for servers $i \in [1, k]$, we require a service characterization for each of the servers that can be composed using the independence assumption. We derive a suitable model using statistical envelope functions of the type of [20].

Definition 2. An arrival process $A(m, n)$ has a statistical sample path envelope ρ_A with error profile $\varepsilon_A(\tau_A)$ for $\tau_A \geq 0$ if it holds for all $n \geq 1$ that

$$\mathbb{P} \left[\max_{\nu \in [1, n]} \{ \rho_A(n - \nu) - A(\nu, n) \} > \tau_A \right] \leq \varepsilon_A(\tau_A).$$

Similarly, a service process $S(m, n)$ has envelope ρ_S with error profile $\varepsilon_S(\tau_S)$ for $\tau_S \geq 0$ if it holds for all $n \geq 1$ that

$$\mathbb{P} \left[\max_{\nu \in [1, n]} \{ S(\nu, n) - \rho_S(n - \nu + 1) \} > \tau_S \right] \leq \varepsilon_S(\tau_S).$$

Corollary 1. Given iid inter-arrival and service times with parameters $\rho_A(-\theta_A)$ (9) and $\rho_S(\theta_S)$ (10) for $\theta_A, \theta_S > 0$. It follows that $\rho_A(-\theta_A)$ and $\rho_S(\theta_S)$ satisfy Def. 2 with error profile $\varepsilon_A(\tau_A) = e^{-\theta_A \tau_A}$ and $\varepsilon_S(\tau_S) = e^{-\theta_S \tau_S}$, respectively.

The proof is a variation of the GI|GI|1 case of Th. 1. It is given in the appendix. We note that a similar envelope characterization that differs mainly by a pre-factor of the error profile can also be obtained for the G|G|1 case, see [23] for a general derivation. We omit the result for reasons of space.

The important quality of the envelopes specified in Def. 2 is the consideration of sample paths of the arrivals and the service by use of the max operator. As a result, the service envelope can be reformulated as

$$\mathbb{P} [\forall \nu \in [1, n] : S(\nu, n) \leq \rho_S(n - \nu + 1) + \tau_S] \geq 1 - \varepsilon_S(\tau_S).$$

By substitution of the envelope for $S(\nu, n)$ into (4) we obtain

$$\mathbb{P} \left[D(n) \leq \max_{\nu \in [1, n]} \{ A(\nu) + \rho_S(n - \nu + 1) \} + \tau_S \right] \geq 1 - \varepsilon_S(\tau_S),$$

so that finally

$$\mathbb{P} \left[D(n) > \max_{\nu \in [1, n]} \{ A(\nu) + \rho_S(n - \nu + 1) \} + \tau_S \right] \leq \varepsilon_S(\tau_S), \quad (25)$$

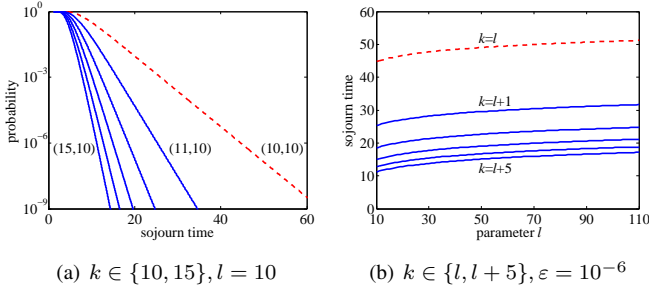


Fig. 8. (k, l) fork-join systems achieve a faster tail decay. Increasing redundancy realizes a diminishing improvement.

for $n \geq 1$. By insertion of (25) into $T(n) = D(n) - A(n)$ a statistical bound of the sojourn time follows for $n \geq 1$ as

$$\mathbb{P}\left[T(n) > \max_{\nu \in [1, n]} \{\rho_S(n - \nu + 1) - A(\nu, n)\} + \tau_S\right] \leq \varepsilon_S(\tau_S).$$

A similar substitution of the arrival envelope from Def. 2 for $A(\nu, n)$ and using the union bound yields

$$\begin{aligned} \mathbb{P}\left[T(n) > \max_{\nu \in [1, n]} \{\rho_S(n - \nu + 1) - \rho_A(n - \nu)\} + \tau_A + \tau_S\right] \\ \leq \varepsilon_A(\tau_A) + \varepsilon_S(\tau_S), \end{aligned}$$

so that for $\rho_S \leq \rho_A$ we have

$$\mathbb{P}[T(n) > \tau_A + \tau_S + \rho_S] \leq \varepsilon_A(\tau_A) + \varepsilon_S(\tau_S). \quad (26)$$

Finally, Cor. 1 gives $\varepsilon_A(\tau_A) = e^{-\theta_A \tau_A}$ and $\varepsilon_S(\tau_S) = e^{-\theta_S \tau_S}$ for any $\theta_A, \theta_S > 0$ that satisfy $\rho_S(\theta_S) \leq \rho_A(-\theta_A)$. We note that (26) does not assume independence of the arrivals and the service.

Considering (k, l) fork-join systems, job $n \geq 1$ departs once l out of its k tasks are finished. The departure time follows as

$$D(n) = \min_{C \in \mathbb{C}(k, l)} \left\{ \max_{i \in C} \{D_i(n)\} \right\}, \quad (27)$$

where $\mathbb{C}(k, l)$ is the set of all combinations of size l that can be chosen from k . For the special case of a standard fork-join system we have $l = k$ so that \mathbb{C} comprises one combination that contains all $i \in [1, k]$ so that (27) becomes $D(n) = \max_{i \in [1, k]} \{D_i(n)\}$. On the other hand, for a system with k redundant jobs we have $l = 1$ and \mathbb{C} comprises k different combinations that each contain one element, so that $D(n) = \min_{i \in [1, k]} \{D_i(n)\}$. In general, (27) considers the l th order statistic of $D(n)$.

Next, we derive a guarantee of the type of (25) for $D(n)$ as defined by (27). Given a service envelope for each of the k parallel servers $S_i(\nu, n)$ as in Def. 2, the departures $D_i(n)$ satisfy (25) for $n \geq 1$ and $i \in [1, k]$. Since (25) takes the form $\mathbb{P}[D_i(n) > x] \leq p_i$ and conversely $\mathbb{P}[D_i(n) \leq x] \geq 1 - p_i$, (25) can be viewed as independent Bernoulli random variables, assuming independence of $S_i(\nu, n)$ for $i \in [1, k]$. In the homogeneous case we have $\rho_{S_i} = \rho_S$ and $p_i = p$ for all

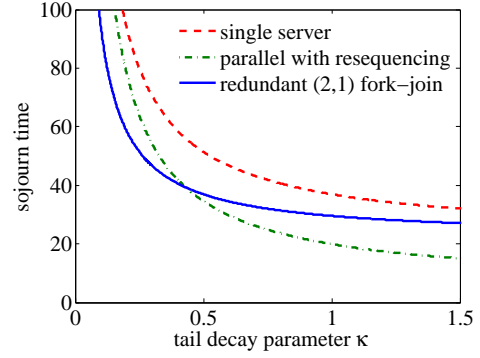


Fig. 9. Two parallel latency-rate servers. The optimal strategy, thinning and resequencing or $(2, 1)$ fork-join, depends on the tail decay parameter κ .

$i \in [1, k]$. Since the number of successful Bernoulli trials is binomial, we obtain from (27) for $n \geq 1$ that

$$\mathbb{P}[D(n) > x] \leq \sum_{j=0}^{l-1} \binom{k}{j} (1-p)^j p^{k-j} =: \varepsilon_{(k, l)}(\tau_S), \quad (28)$$

where $x = \max_{\nu \in [1, n]} \{A(\nu) + \rho_S(n - \nu + 1)\} + \tau_S$ and $p = e^{-\theta_S \tau_S}$ from Cor. 1. A statistical bound of the sojourn time follows from (26) by substitution of $\varepsilon_{(k, l)}(\tau_S)$ for $\varepsilon_S(\tau_S)$.

D|M|I Servers: To focus only on the effect of (k, l) fork-join systems, we use deterministic inter-arrival times ρ_A , i.e., $A(\nu, n) = \rho_A(n - \nu)$. The service times are iid exponential with parameters $\mu_i = \mu$ for all servers $i \in [1, k]$ so that $\rho_S(\theta_S) = \ln(\mu/(\mu - \theta_S))/\theta_S$ for $\theta_S \in (0, \mu)$ from (13). We choose the largest θ_S that satisfies $\rho_S(\theta_S) \leq \rho_A$. In Fig. 8, we show the tail decay of the sojourn time (26) where the (k, l) fork-join system is governed by (28). The parameters are $\rho_A = 1.25$ and $\mu = 1$. Fig. 8 shows how redundancy, i.e., $k > l$, improves the sojourn time. It becomes apparent that few redundant servers can achieve a significant advantage. Adding more redundancy realizes a diminishing improvement.

Latency-rate Servers: The system model (25) takes the basic form of a latency-rate function. It can be parameterized by $1/\rho_S$ that has the interpretation of a service rate, and τ_S that is a latency with error profile $\varepsilon_S(\tau_S)$. We consider $k = 2$ parallel servers, each with $\rho_S = 1$ and $\varepsilon_S(\tau_S) = e^{-\kappa \tau_S}$ where $\kappa > 0$ is the speed of the tail decay that determines the average latency $1/\kappa$. An example is a packet data stream that can be routed via two different paths. The arrivals have exponential inter-arrival times with parameter $\lambda = 0.7$. In Fig. 9, we show bounds of the sojourn time with $\varepsilon = 10^{-6}$, where we use κ as a parameter to evaluate three fundamental strategies: 1) use of only one server; 2) use of both servers in parallel with deterministic thinning and resequencing; 3) use of both servers as a $(2, 1)$ fork-join system, i.e., with 1-out-of-2 redundancy.

If κ is small, the servers frequently experience large delays. Since these delays occur independently for each of the two servers, the redundancy of the $(2, 1)$ fork-join system provides a significant advantage. In contrast, if κ is large, delays are mostly due to the variability of the arrivals. In case of the $(2, 1)$ fork-join system, this affects both servers in the same way, so

that the redundant server provides little utility. In this case, the thinning strategy is clearly superior as it benefits from load balancing compared to the two other strategies. Thinning can, however, only combat delays that are due to the variability and the load caused by the arrivals. In contrast, in case of small κ , thinning and resequencing can result in significant additional delays, as jobs have to wait more frequently at the resequencing stage to ensure in-order delivery. This effect will eventually, if κ is reduced further, result in increased sojourn times compared to the single server system.

B. Multi-Stage Fork-Join Networks

We derive end-to-end results for multi-stage fork-join networks, where we consider h fork-join stages in tandem, each with k parallel servers. We use subscript $i \in [1, k]$ to distinguish the servers of a stage and superscript $j \in [1, h]$ to denote the stages. We employ the model from Sec. IV-A and specify each server by (25) with parameters ρ_{S_i} and $\varepsilon_{S_i}(\tau_S)$. For notational convenience, we consider homogeneous servers, i.e., the service processes at each of the servers have the same distribution so that we can drop the subscript i . It follows that each stage also satisfies (25) with parameters $\rho_{stg} = \rho_S$ and $\varepsilon_{stg}(\tau_S)$. In general, we have $\varepsilon_{stg}(\tau_S) = k\varepsilon_S(\tau_S)$ from the union bound and for (k, l) fork-join systems with independent parallel servers $\varepsilon_{stg}(\tau_S) = \varepsilon_{(k,l)}(\tau_S)$ from (28).

In a tandem of fork-join stages, the departures of stage j become the arrivals of stage $j+1$, i.e., $A^{j+1}(n) = D^j(n)$ for $j \in [1, h-1]$. As $D^j(n)$ is given by (25) for all $j \in [1, h]$, it would appear that a solution can be obtained by recursive insertion of (25) for each stage. The difficulty is, however, that (25) evaluates sample paths of $A^j(n)$ but does not provide a sample path guarantee for $D^j(n)$, see [28] for a discussion of the challenges. To evaluate sample paths of $D^j(n)$ we adapt the basic method from [29] to max-plus systems and define a sample path guarantee for $D^j(n)$ for any $m \geq 1$ as

$$\mathbb{P} \left[\exists n \in [1, m] : D^j(n) > \max_{\nu \in [1, n]} \{A^j(\nu) + \rho_S(n - \nu + 1)\} + \delta(m - n + 1) + \tau_S \right] \leq \varepsilon_{stg}^\delta(\tau_S), \quad (29)$$

where $\delta > 0$ is a free parameter. Given (25) holds for stage j with $\varepsilon_{stg}(\tau_S)$ as defined above, we derive by the union bound that

$$\varepsilon_{stg}^\delta(\tau_S) = \sum_{n=1}^m \varepsilon_{stg}(\tau_S + \delta(m - n + 1)). \quad (30)$$

Using (29), a characterization of the end-to-end service of a network of h stages can be derived recursively, starting from the last stage. For stage h , we obtain from (29)¹ for $n = m$

¹Alternatively, for the special case of stage h , (25) can be used to obtain a tighter guarantee. We use (29) here for notational simplicity.

and $m \geq 1$ that

$$\mathbb{P} \left[D^h(m) > \max_{\nu \in [1, m]} \{A^h(\nu) + \rho_S(m - \nu + 1)\} + \delta + \tau_S \right] \leq \varepsilon_{stg}^\delta(\tau_S).$$

Since $A^h(\nu) = D^{h-1}(\nu)$ and $\nu \in [1, m]$ we can use (29) again, to derive

$$\mathbb{P} \left[D^h(m) > \max_{\nu \in [1, m]} \left\{ \max_{\vartheta \in [1, \nu]} \{A^{h-1}(\vartheta) + \rho_S(\nu - \vartheta + 1)\} + \delta(m - \nu + 1) + \tau_S + \rho_S(m - \nu + 1) \right\} + \delta + \tau_S \right] \leq 2\varepsilon_{stg}^\delta(\tau_S),$$

where we used the union bound to estimate the probability. After some reordering we estimate $\delta(m - \nu + 1) \leq \delta(m - \vartheta + 1)$, since $\vartheta \in [1, \nu]$. Finally, we can combine the outer and the inner max to arrive at

$$\mathbb{P} \left[D^h(m) > \max_{\vartheta \in [1, m]} \{A^{h-1}(\vartheta) + \rho_S(m - \vartheta + 2)\} + \delta(m - \vartheta + 1) + \delta + 2\tau_S \right] \leq 2\varepsilon_{stg}^\delta(\tau_S).$$

Repeating the same steps, we obtain by recursive insertion of (29) for $m \geq 1$ that

$$\mathbb{P} \left[D^h(m) > \max_{\nu \in [1, m]} \{A^1(\nu) + \rho_S(m - \nu + h)\} + (h - 1)\delta(m - \nu + 1) + \delta + h\tau_S \right] \leq h\varepsilon_{stg}^\delta(\tau_S). \quad (31)$$

The end-to-end sojourn time of the network follows as before by insertion of (31) into $T^{e2e}(m) = D^h(m) - A^1(m)$.

GI|GI|I servers: We derive $\varepsilon_{stg}^\delta(\tau_S)$ for a fork-join system with k servers each with iid service times to show how bounds of $T^{e2e}(m)$ scale with k and h . With Cor. 1 we have $\varepsilon_{stg}(\tau_S) = ke^{-\theta_S \tau_S}$ so that (30) can be estimated as

$$\begin{aligned} \varepsilon_{stg}^\delta(\tau_S) &= \sum_{n=1}^m ke^{-\theta_S(\tau_S + \delta(m - n + 1))} \\ &\leq ke^{-\theta_S \tau_S} \sum_{n=1}^m e^{-\theta_S \delta n} \\ &\leq ke^{-\theta_S \tau_S} \int_0^\infty e^{-\theta_S \delta y} dy \\ &= \frac{ke^{-\theta_S \tau_S}}{\theta_S \delta}, \end{aligned} \quad (32)$$

where we used that $e^{-\theta_S \delta n} \leq \int_{n-1}^n e^{-\theta_S \delta y} dy$ for $n \geq 1$ since $e^{-\theta_S \delta y}$ is a decreasing function in $y \geq 0$. Next, we define a constant $\beta > 0$ where we let $\delta = \beta/h$ and insert (32) into (31). It follows that

$$h\varepsilon_{stg}^\delta(\tau_S) = e^{-\theta_S(\tau_S - \frac{1}{\theta_S} \ln \frac{h^2 k}{\theta_S \beta})} = e^{-\theta_S z},$$

where we define z so that $\tau_S = z + \frac{1}{\theta_S} \ln \frac{h^2 k}{\theta_S \beta}$. By insertion of this expression into (31) and using another variable substitution, we arrive at

$$\mathbb{P} \left[D^h(m) > \max_{\nu \in [1, m]} \{A^1(\nu) + (\rho_S(\theta_S) + \beta)(m - \nu)\} + \beta + h \left(\rho_S(\theta_S) + \tau_S + \frac{1}{\theta_S} \ln \frac{h^2 k}{\theta_S \beta} \right) \right] \leq e^{-\theta_S \tau_S}. \quad (33)$$

A statistical upper bound of the end-to-end sojourn time $T^{e2e}(m) = D^h(m) - A^1(m)$ follows with (33) as

$$\mathbb{P} \left[T^{e2e}(m) > \max_{\nu \in [1, m]} \{(\rho_S(\theta_S) + \beta)(m - \nu) - A^1(\nu, m)\} + \beta + h \left(\rho_S(\theta_S) + \tau_S + \frac{1}{\theta_S} \ln \frac{h^2 k}{\theta_S \beta} \right) \right] \leq e^{-\theta_S \tau_S}.$$

Considering arrivals with envelope $\rho_A(-\theta_A)$ as in Def. 2 and Cor. 1, it follows for $\rho_S(\theta_S) + \beta \leq \rho_A(-\theta_A)$ that

$$\mathbb{P} \left[T^{e2e}(m) > \beta + \tau_A + h \left(\rho_S(\theta_S) + \tau_S + \frac{1}{\theta_S} \ln \frac{h^2 k}{\theta_S \beta} \right) \right] \leq e^{-\theta_A \tau_A} + e^{-\theta_S \tau_S}.$$

which grows with $\mathcal{O}(h \ln(hk))$. The result extends the fundamental scaling $\mathcal{O}(h \ln h)$ for h tandem servers from [29] to h tandem fork-join stages each with k parallel servers. It holds without assuming independent servers within or across stages.

V. CONCLUSIONS

We presented a general model of fork-join systems in max-plus system theory and derived essential solutions for k parallel $G|G|1$ servers, that generalize previous $\ln k$ scaling results. Beyond standard fork-join systems, we considered parallel servers with resequencing, where each server has a thinned arrival process. We included applications to load balancing that showed the significant potential but also revealed fundamental limitations. Our results showed that deterministic thinning is superior to random thinning. Taking advantage of an envelope-based approach, we contributed solutions for advanced (k, l) fork-join systems and demonstrated the substantial benefit of using independent redundant servers. We provided insights into the configuration of redundant servers versus load-balancing. We extended the results to multi-stage fork-join networks with h stages, where we discovered a scaling with $\mathcal{O}(h \ln(hk))$.

VI. APPENDIX

Proof of Th. 1: We only show the proof of the sojourn time, as the proof of the waiting time follows similarly.

$G|G|1$ servers: We obtain from (5) for $\theta > 0$ that

$$\mathbb{E}[e^{\theta T(n)}] \leq \sum_{\nu=1}^n \mathbb{E}[e^{\theta S(\nu, n)}] \mathbb{E}[e^{-\theta A(\nu, n)}],$$

where we estimated the maximum by the sum of its arguments and used the statistical independence of arrivals and service.

By insertion of the (σ, ρ) -constraints from Def. 1 we have

$$\mathbb{E}[e^{\theta T(n)}] \leq e^{\theta(\sigma_A(-\theta) + \sigma_S(\theta) + \rho_S(\theta))} \sum_{\nu=1}^n e^{-\theta(\rho_A(-\theta) - \rho_S(\theta))(n - \nu)}.$$

Next, we estimate

$$\begin{aligned} \sum_{\nu=1}^n e^{-\theta(\rho_A(-\theta) - \rho_S(\theta))(n - \nu)} &\leq \sum_{\nu=0}^{\infty} (e^{-\theta(\rho_A(-\theta) - \rho_S(\theta))})^\nu \\ &= \frac{1}{1 - e^{-\theta(\rho_A(-\theta) - \rho_S(\theta))}}, \end{aligned}$$

where we used the geometric sum for $\rho_S(\theta) < \rho_A(-\theta)$. By use of Chernoff's bound $\mathbb{P}[X \geq x] \leq e^{-\theta x} \mathbb{E}[e^{\theta X}]$ we obtain

$$\mathbb{P}[T(n) \geq \tau] \leq \frac{e^{\theta(\sigma_A(-\theta) + \sigma_S(\theta))}}{1 - e^{-\theta(\rho_A(-\theta) - \rho_S(\theta))}} e^{\theta \rho_S(\theta)} e^{-\theta \tau}.$$

$GI|GI|1$ servers: From (5) we have

$$T(n) = \max_{m \in [1, n]} \{S(n - m + 1, n) - A(n - m + 1, n)\}.$$

For $\theta > 0$ we can write

$$\begin{aligned} \mathbb{P}[T(n) > \tau] &= \mathbb{P} \left[\max_{m \in [1, n]} \left\{ e^{\theta(S(n - m + 1, n) - A(n - m + 1, n))} \right\} > e^{\theta \tau} \right]. \end{aligned}$$

Now consider the process

$$U(m) = e^{\theta(S(n - m + 1, n) - A(n - m + 1, n))}.$$

By definition of $A(m, n)$ (1) and $S(m, n)$ (2) it follows that

$$U(m + 1) = U(m) e^{\theta(S(n - m) - A(n - m, n - m + 1))}.$$

The conditional expectation can be computed as

$$\begin{aligned} \mathbb{E}[U(m + 1) | U(m), U(m - 1), \dots, U(1)] \\ = U(m) \mathbb{E}[e^{\theta S(n - m)}] \mathbb{E}[e^{-\theta A(n - m, n - m + 1)}], \end{aligned}$$

where we used the statistical independence of the inter-arrival and service times. If $\rho_S(\theta) \leq \rho_A(-\theta)$, it holds that $\mathbb{E}[e^{\theta S(n - m)}] \mathbb{E}[e^{-\theta A(n - m, n - m + 1)}] \leq 1$ and

$$\mathbb{E}[U(m + 1) | U(m), U(m - 1), \dots, U(1)] \leq U(m),$$

i.e., $U(m)$ is a supermartingale. By application of Doob's inequality for submartingales [36, Theorem 3.2, p. 314] and the formulation for supermartingales [24], [38] we have for non-negative $U(m)$ for $m \geq 1$ that

$$x \mathbb{P} \left[\max_{m \in [1, n]} \{U(m)\} \geq x \right] \leq \mathbb{E}[U(1)]. \quad (34)$$

We derive

$$\mathbb{E}[U(1)] = \mathbb{E}[e^{\theta(S(n, n) - A(n, n))}] = \mathbb{E}[e^{\theta S(1)}].$$

Letting $x = e^{\theta \tau}$ we have from (34) that

$$\mathbb{P}[T(n) \geq \tau] \leq e^{\theta \rho_S(\theta)} e^{-\theta \tau}.$$

Finally, using the union bound to estimate (7) gives Th. 1. ■

Proof of Cor. 1: We only show the proof of the arrival envelope as the proof of the service envelope follows similarly.

$GI|GI|I$ servers: For $\theta > 0$ we can rewrite Def. 2 as

$$\begin{aligned} & \mathbb{P} \left[\max_{\nu \in [1, n]} \{ \rho_A(n - \nu) - A(\nu, n) \} > \tau_A \right] \\ &= \mathbb{P} \left[\max_{m \in [1, n]} \{ e^{\theta(\rho_A(m-1) - A(n-m+1, n))} \} > e^{\theta\tau_A} \right]. \end{aligned}$$

Now consider the process

$$U(m) = e^{\theta(\rho_A(m-1) - A(n-m+1, n))}.$$

By definition of $A(m, n)$ (1) it follows that

$$U(m+1) = U(m)e^{\theta(\rho_A - A(n-m, n-m+1))},$$

which has conditional expectation

$$\begin{aligned} & \mathbb{E}[U(m+1)|U(m), U(m-1), \dots, U(1)] \\ &= U(m)e^{\theta\rho_A} \mathbb{E}[e^{-\theta A(n-m, n-m+1)}]. \end{aligned}$$

By substitution of $\rho_A = \rho_A(-\theta)$ from (9) it follows that $\mathbb{E}[U(m+1)|U(m), U(m-1), \dots, U(1)] = U(m)$ is a martingale and hence also a supermartingale. With (34), $x = e^{\theta\tau_A}$, and since $\mathbb{E}[U(1)] = 1$ it holds that

$$\mathbb{P} \left[\max_{\nu \in [1, n]} \{ \rho_A(n - \nu) - A(\nu, n) \} \geq \tau_A \right] \leq e^{-\theta\tau_A},$$

which proves Cor. 1. \blacksquare

REFERENCES

- [1] J. Tan, X. Meng, and L. Zhang, "Delay tails in MapReduce scheduling," *ACM Sigmetrics Perf. Eval. Rev.*, vol. 40, no. 1, pp. 5–16, Jun. 2012.
- [2] F. Baccelli, A. M. Makowski, and A. Schwartz, "The fork-join queue and related systems with synchronization constraints: Stochastic ordering and computable bounds," *Adv. in Appl. Probab.*, vol. 21, no. 3, pp. 629–660, Sep. 1989.
- [3] B. Kemper and M. Mandjes, "Mean sojourn times in two-queue fork-join systems: bounds and approximations," *OR Spektrum*, vol. 34, no. 3, pp. 723–742, Jul. 2012.
- [4] L. Flatto and S. Hahn, "Two parallel queues created by arrivals with two demands," *SIAM J. Appl. Math.*, vol. 44, no. 5, pp. 1041–1053, 1984.
- [5] R. Nelson and A. N. Tantawi, "Approximate analysis of fork/join synchronization in parallel queues," *IEEE Trans. Comput.*, vol. 37, no. 6, pp. 739–743, Jun. 1988.
- [6] A. S. Lebrecht and W. J. Knottenbelt, "Response time approximations in fork-join queues," in *Proc. of UKPEW*, Jul. 2007.
- [7] S.-S. Ko and R. F. Serfozo, "Sojourn times in G/M/1 fork-join networks," *Naval Research Logistics*, vol. 55, no. 5, pp. 432–443, May 2008.
- [8] X. Tan and C. Knessl, "A fork-join queueing model: Diffusion approximation, integral representations and asymptotics," *Queueing Systems*, vol. 22, no. 3-5, pp. 287–322, Sep. 1996.
- [9] S. Varma and A. M. Makowski, "Interpolation approximations for symmetric fork-join queues," *Performance Evaluation*, vol. 20, no. 1-3, pp. 245–265, May 1994.
- [10] E. Varki, "Mean value technique for closed fork-join networks," *ACM Sigmetrics Perf. Eval. Rev.*, vol. 27, no. 1, pp. 103–112, May 1999.
- [11] F. Alomari and D. A. Menascé, "Efficient response time approximations for multiclass fork and join queues in open and closed queueing networks," *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 6, pp. 1437–1446, Jun. 2014.
- [12] A. Rizk, F. Poloczek, and F. Ciucu, "Computable bounds in fork-join queueing systems," in *Proc. of ACM Sigmetrics*, Jun. 2015.
- [13] G. Kesidis, B. Urgaonkar, Y. Shan, S. Kamarava, and J. Liebeherr, "Network calculus for parallel processing," in *Proc. of MAMA Workshop at ACM SIGMETRICS*, Jun. 2015.
- [14] Y. Han and A. M. Makowski, "Resequencing delays under multipath routing – asymptotics in a simple queueing model," in *Proc. of IEEE INFOCOM*, Apr. 2006.
- [15] Y. Xia and D. Tse, "On the large deviation of resequencing queue size: 2-M/M/1 case," *IEEE Trans. Inf. Theory*, vol. 54, no. 9, pp. 4107–4118, Sep. 2008.
- [16] Y. Gao and Y. Q. Zhao, "Large deviations for re-sequencing buffer size," *IEEE Trans. Inf. Theory*, vol. 58, no. 2, pp. 1003–1009, Feb. 2012.
- [17] P. Harrison and S. Zertal, "Queueing models with maxima of service times," in *Proc. of TOOLS*, Sep. 2003, pp. 152–168.
- [18] G. Joshi, Y. Liu, and E. Soljanin, "On the delay-storage trade-off in content download from coded distributed storage systems," *IEEE J. Sel. Areas Commun.*, vol. 32, no. 5, pp. 989–997, May 2014.
- [19] N. B. Shah, K. Lee, and K. Ramchandran, "The MDS queue: Analyzing the latency performance of erasure codes," in *Proc. of ISIT*, Jul. 2014.
- [20] Q. Yin, Y. Jiang, S. Jiang, and P. Y. Kong, "Analysis of generalized stochastically bounded bursty traffic for communication networks," in *Proc. of IEEE LCN*, Nov. 2002, pp. 141–149.
- [21] F. Baccelli, G. Cohen, G. J. Olsder, and J.-P. Quadrat, *Synchronization and Linearity: An Algebra for Discrete Event Systems*. Wiley, 1992.
- [22] C.-S. Chang, *Performance Guarantees in Communication Networks*. Springer-Verlag, 2000.
- [23] J. Xie and Y. Jiang, "Stochastic network calculus models under max-plus algebra," in *Proc. of IEEE GLOBECOM*, 2009.
- [24] Y. Jiang, "Network calculus and queueing theory: Two sides of one coin," in *Proc. of VALUETOOLS*, 2009, pp. 1–12, Invited Paper.
- [25] R. Lübben, M. Fidler, and J. Liebeherr, "Stochastic bandwidth estimation in networks with random service," *IEEE/ACM Trans. Netw.*, vol. 22, no. 2, pp. 484–497, Apr. 2014.
- [26] R. L. Cruz, "A calculus for network delay, part I and II: Network elements in isolation and network analysis," *IEEE Trans. Inf. Theory*, vol. 37, no. 1, pp. 114–141, Jan. 1991.
- [27] J.-Y. Le Boudec and P. Thiran, *Network Calculus A Theory of Deterministic Queueing Systems for the Internet*. Springer-Verlag, 2001.
- [28] A. Burchard, J. Liebeherr, and S. Patek, "A min-plus calculus for end-to-end statistical service guarantees," *IEEE Trans. Inf. Theory*, vol. 52, no. 9, pp. 4105–4114, Aug. 2006.
- [29] F. Ciucu, A. Burchard, and J. Liebeherr, "Scaling properties of statistical end-to-end bounds in the network calculus," *IEEE/ACM Trans. Netw.*, vol. 14, no. 6, pp. 2300–2312, Jun. 2006.
- [30] M. Fidler, "An end-to-end probabilistic network calculus with moment generating functions," in *Proc. of IWQoS*, Jun. 2006, pp. 261–270.
- [31] Y. Jiang and Y. Liu, *Stochastic Network Calculus*. Springer-Verlag, Sep. 2008.
- [32] M. Fidler, "A survey of deterministic and stochastic service curve models in the network calculus," *IEEE Commun. Surveys Tuts.*, vol. 12, no. 1, pp. 59–86, 2010.
- [33] F. Ciucu and J. Schmitt, "Perspectives on network calculus - no free lunch but still good value," in *Proc. of ACM SIGCOMM*, Aug. 2012, pp. 311–322.
- [34] M. Fidler and A. Rizk, "A guide to the stochastic network calculus," *IEEE Commun. Surveys Tuts.*, vol. 17, no. 1, pp. 92–105, 2015.
- [35] J. F. C. Kingman, "A martingale inequality in the theory of queues," *Math. Proc. Cambridge*, vol. 60, no. 2, pp. 359–361, Apr. 1964.
- [36] J. L. Doob, *Stochastic Processes*. Wiley, 1953.
- [37] I. Adan and J. Resing, *Queueing Systems*. Eindhoven University of Technology, Department of Mathematics and Computing Science, 2015.
- [38] Y. Jiang, "A note on applying stochastic network calculus," Tech. Rep., 2010.